

Online Speculative Decoding

– Towards accelerated LLM inference in the wild

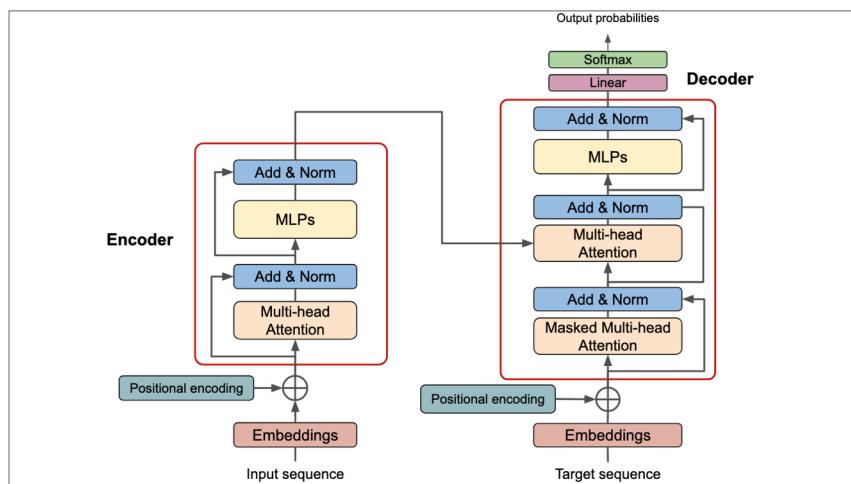
Zhijie Deng

* Joint work with Xiaoxuan Liu, Lanxiang Hu, Peter Bailis, Ion Stoica, Alvin Cheung, and Hao Zhang

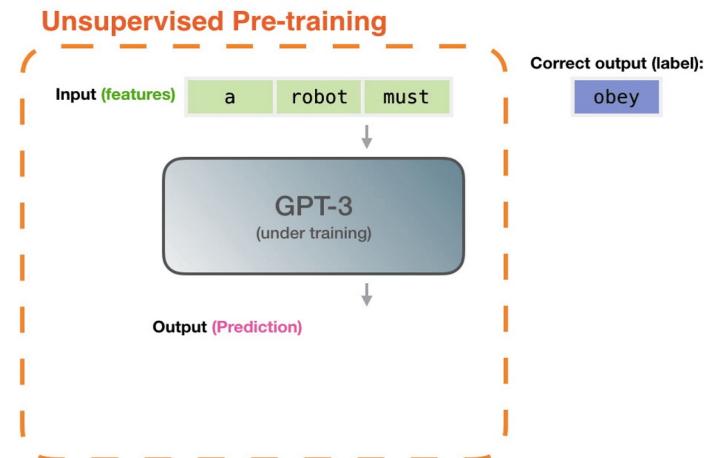
SJTU

<https://arxiv.org/abs/2310.07177>

Large Language Models (LLMs)

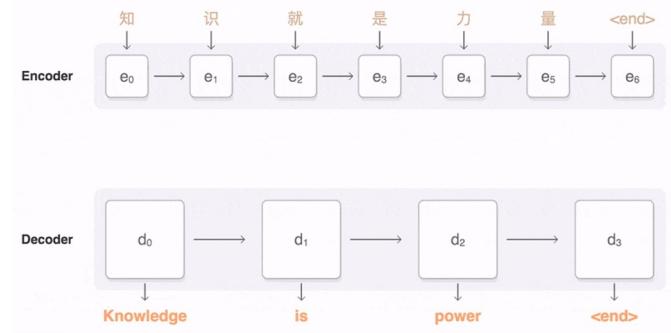


Transformer architecture

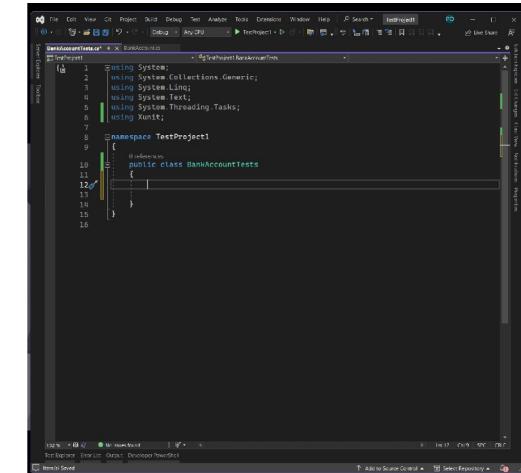


Self-supervised pre-training; 175B parameters; 300B tokens

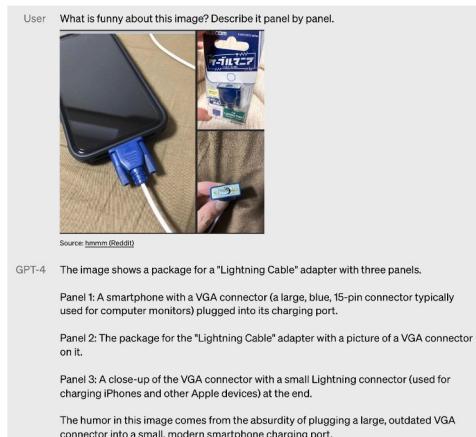
Typical applications of LLMs



Machine translation



Code generation



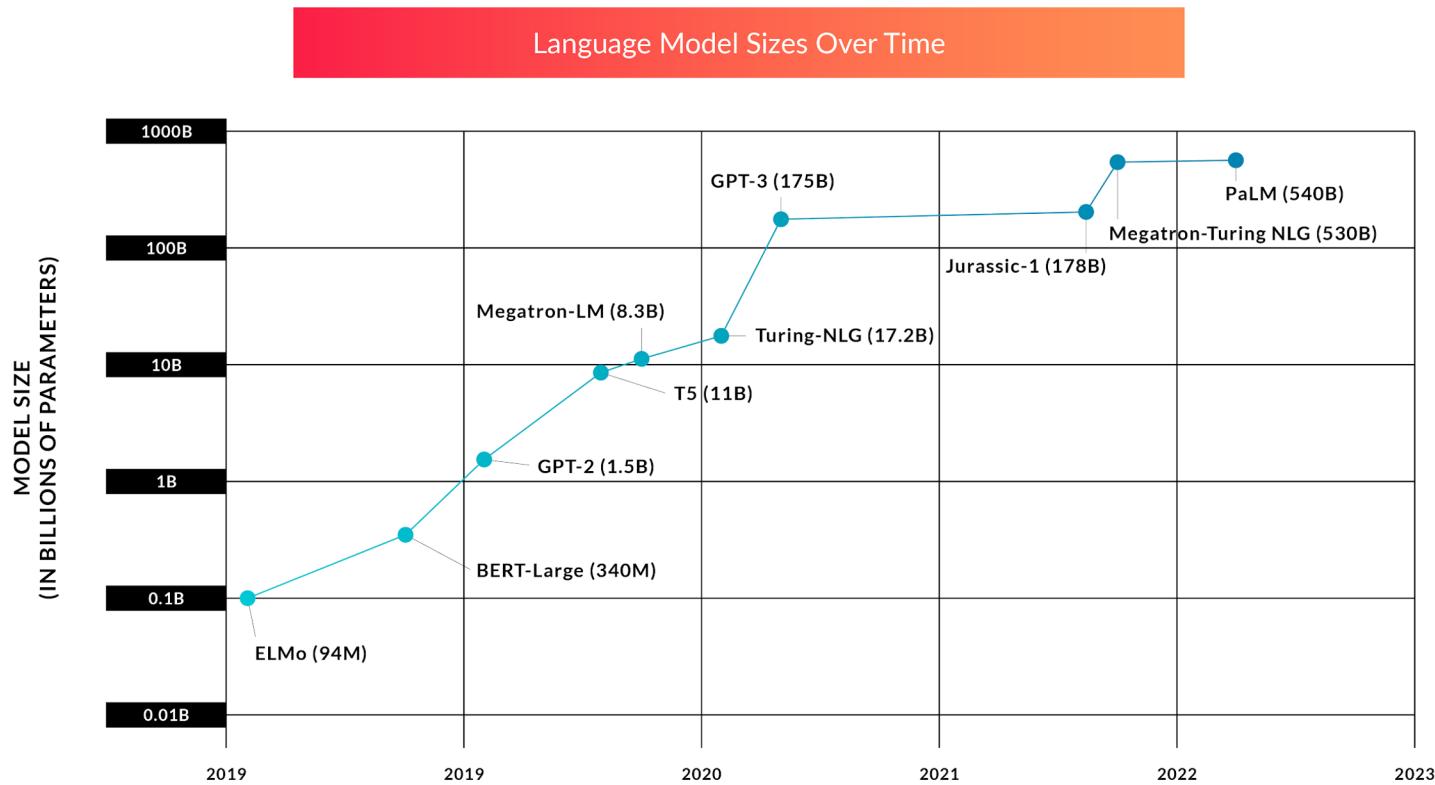
(Multimodal) Chatbot



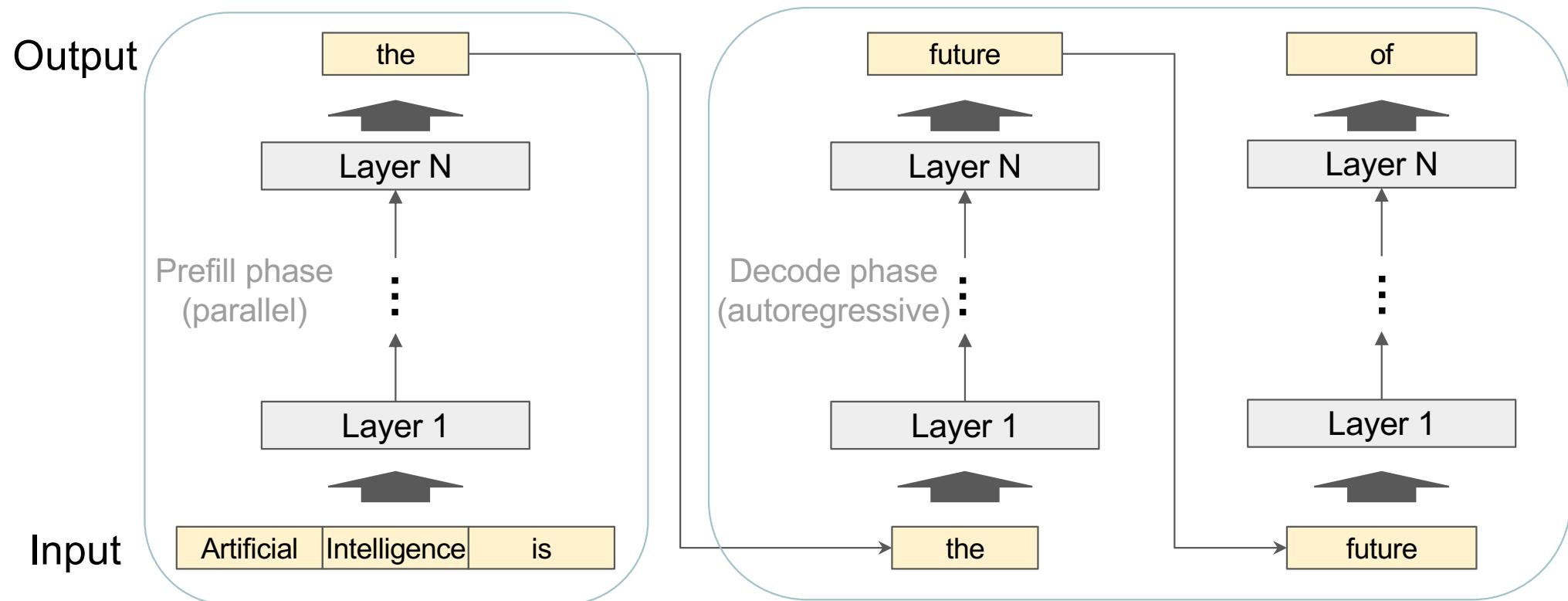
Agents

etc.

LLM inference latency affects the quality of service and user experience



How do LLMs infer?

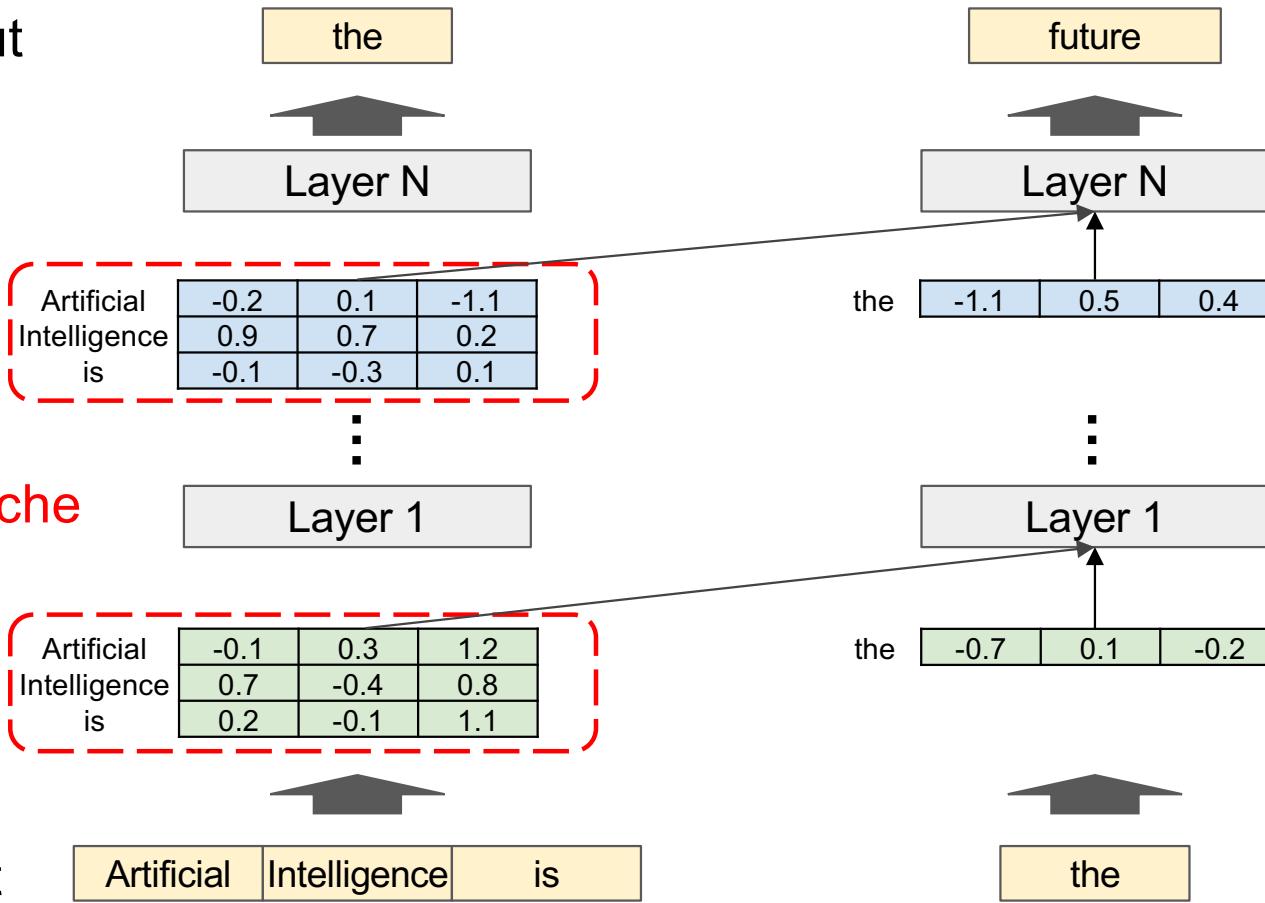


Repeat until the sequence

- Reaches its pre-defined maximum length (e.g., 2048 tokens)
- Generates certain tokens (e.g., "<|end of sequence|>")

KV Cache

Output



KV Cache

Output

KV Cache

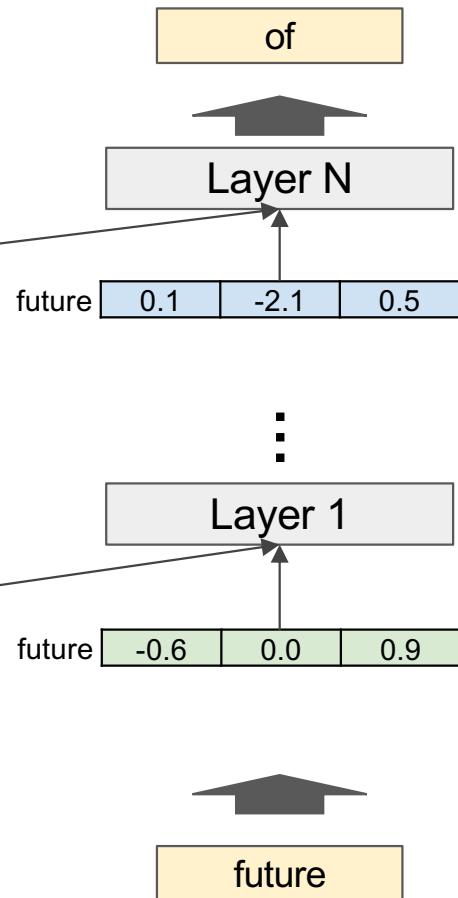
Input

Artificial
Intelligence
is
the

-0.2	0.1	-1.1
0.9	0.7	0.2
-0.1	-0.3	0.1
-1.1	0.5	0.4

Artificial
Intelligence
is
the

-0.1	0.3	1.2
0.7	-0.4	0.8
0.2	-0.1	1.1
-0.7	0.1	-0.2



Two features

- Prefill phase takes about as much time as the generation of each subsequent token due to the use of GPU's parallel compute
- LLM inference is **memory-IO bound** (between High Bandwidth Memory (HBM) and Static Random Access Memory (SRAM)), not compute bound
 - Abundant “**spare FLOPs**” exists in the serving of LLMs

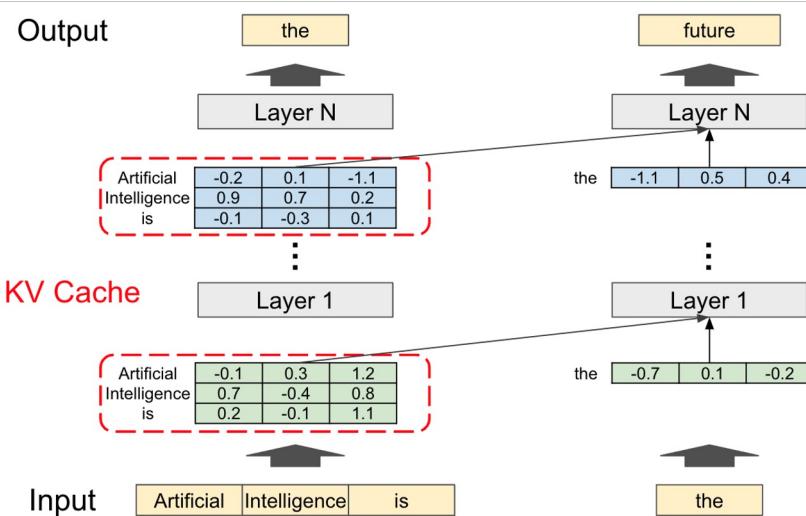
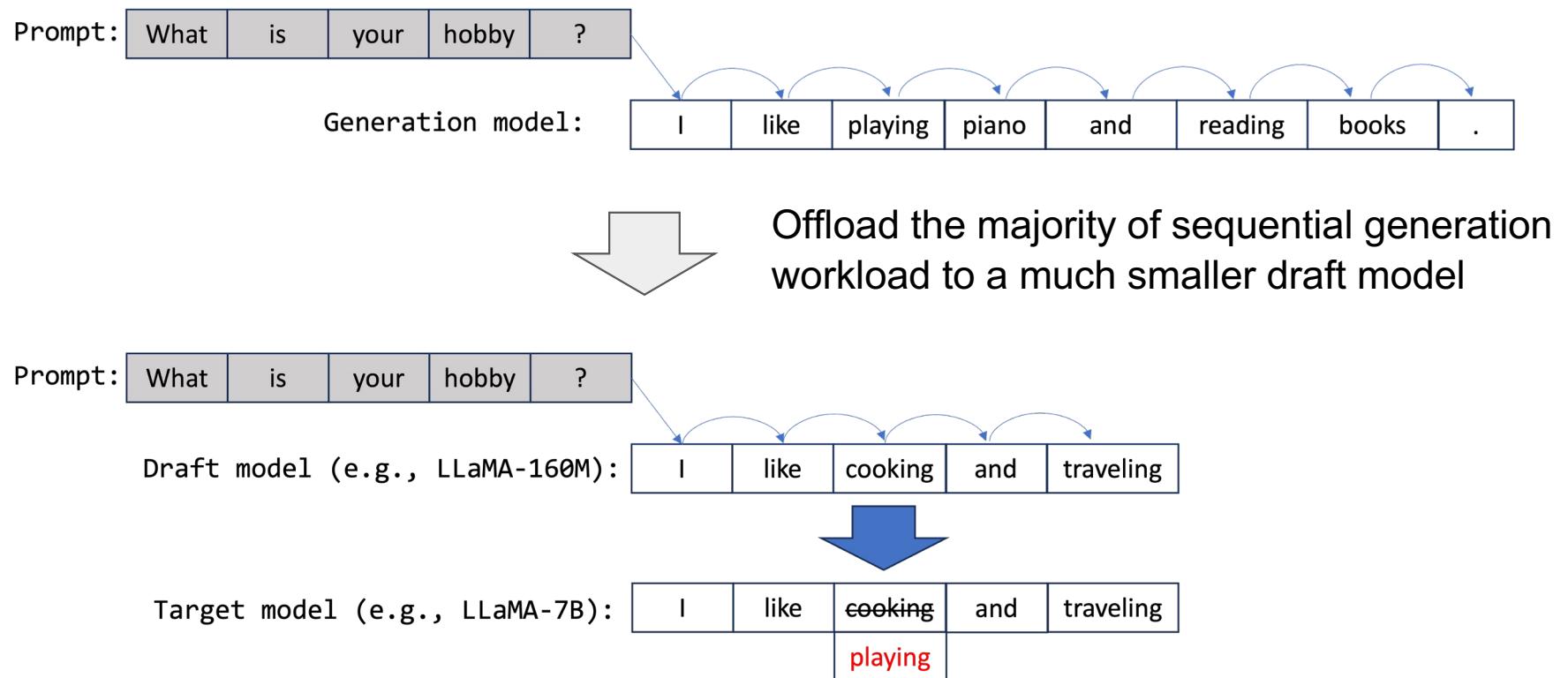


Table 1. Examples of neural network operations with their arithmetic intensities. Limiters assume FP16 data and an NVIDIA V100 GPU.

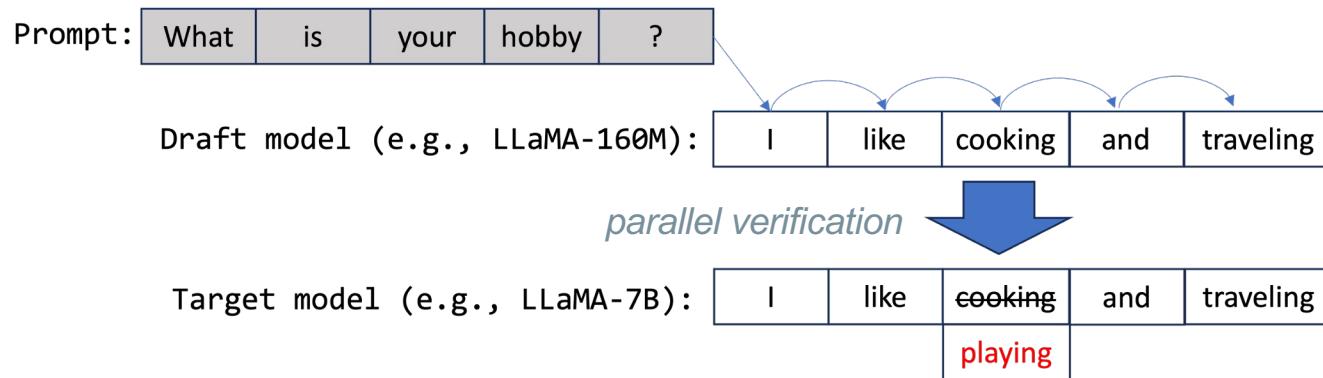
Operation	Arithmetic Intensity	Usually limited by...
Linear layer (4096 outputs, 1024 inputs, batch size 512)	315 FLOPS/B	arithmetic
Linear layer (4096 outputs, 1024 inputs, batch size 1)	1 FLOPS/B	memory
Max pooling with 3x3 window and unit stride	2.25 FLOPS/B	memory
ReLU activation	0.25 FLOPS/B	memory
Layer normalization	< 10 FLOPS/B	memory

<https://docs.nvidia.com/deeplearning/performance/dl-performance-gpu-background/index.html#understand-perf>

Speculative decoding reduces LLM inference latency



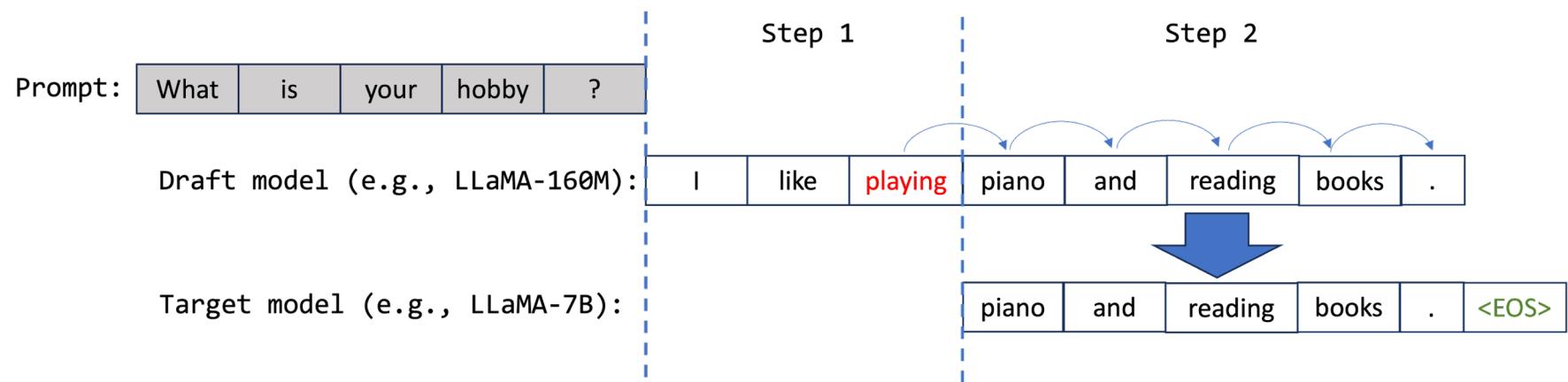
Speculative decoding



Key insights:

- Some tokens are **straightforward** to generate, while others are more **challenging**
- We can utilize a streamlined 'draft' model for the easier tokens
- To ensure identical output to the original generation method, the tokens proposed by draft model are then validated by the target model **in parallel** by the principle of **rejection sampling**

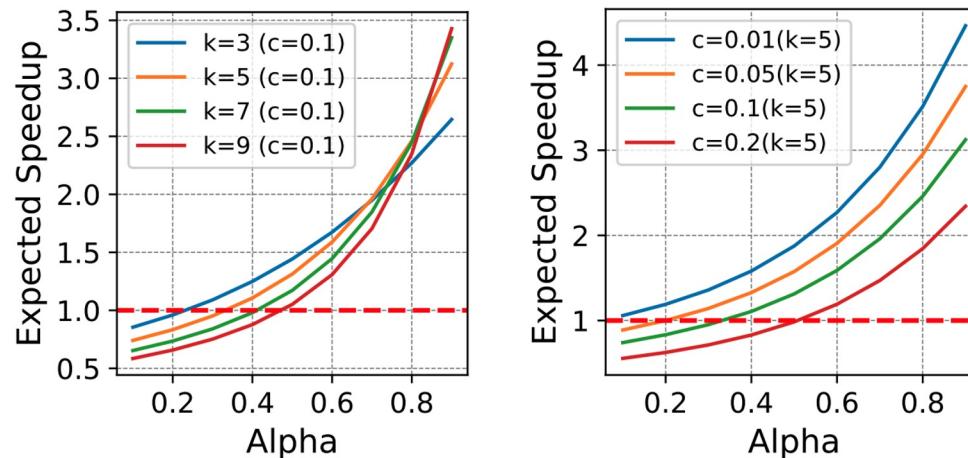
Speculative decoding (cont.)



Why can speculative decoding reduce latency?

- The draft model retains the autoregressive nature, generating tokens one at a time but with a significantly **faster** speed
- The target model can validate multiple generated tokens from the draft model in a **single** forward pass
- Consequently, speculative decoding helps **amortize** the overhead of loading model weights and key-value caches
 - *Originally, each token required accessing the weights and key-value cache, but now it's reduced to just one access per x tokens, where x represents the number of accepted tokens in each generation step*

When will speculative decoding work?

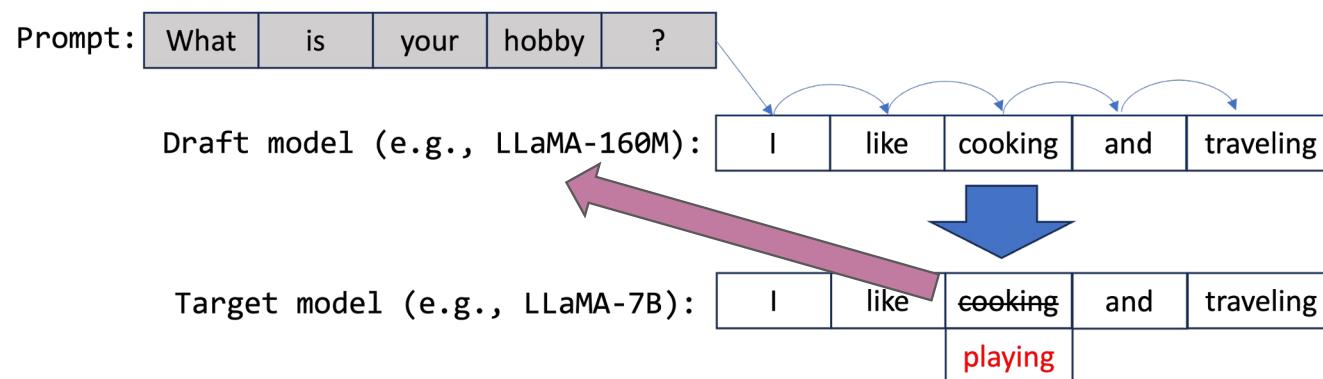


c : the time ratio for a single run between the draft and target model. k : number of proposed tokens each step. Alpha: token acceptance rate.

- **Better token acceptance rate leads to more speedup:** the draft model must approximate the target model sufficiently while being small to achieve latency reduction

Opportunity for improving token acceptance rate online

- Speculative decoding detects inaccuracies within the smaller draft model and provides corrections
 - Such information can be harnessed to refine the draft model, thereby enhancing the draft model's token acceptance rate, all *without incurring any additional labeling costs*

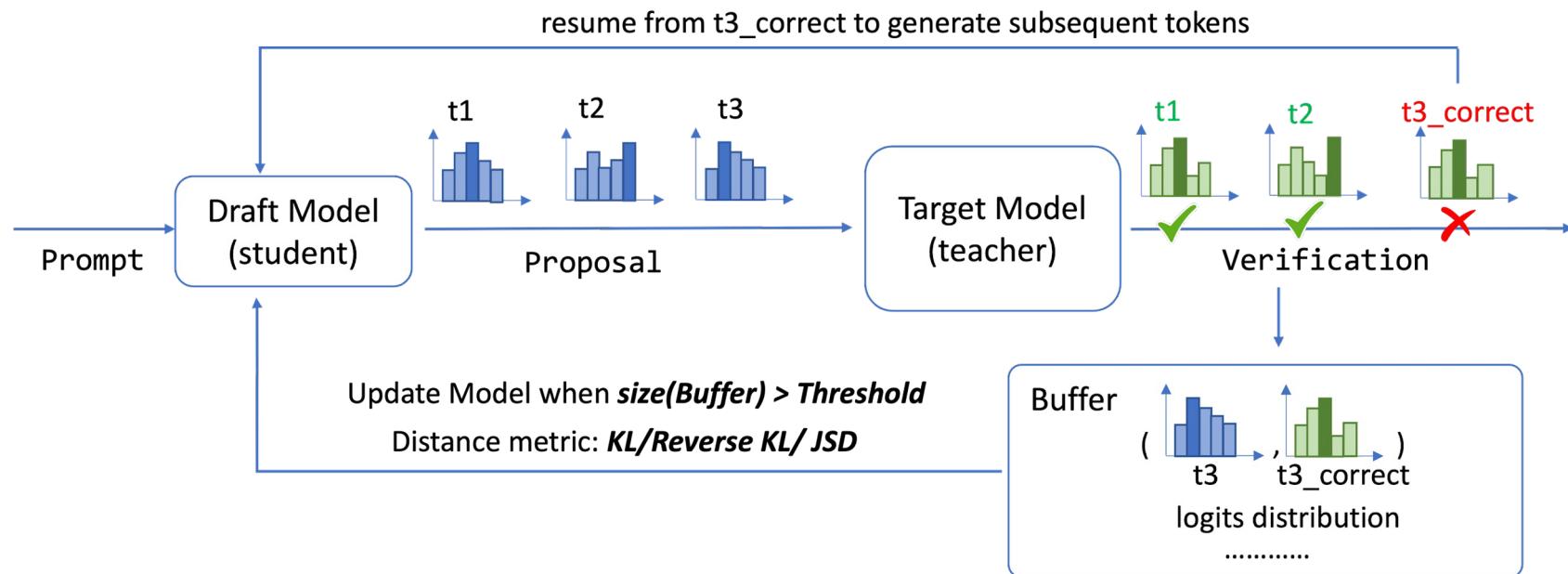


- The *spare FLOPs*

Necessity for improving token acceptance rate online

- Performance of speculative decoding algorithm depends heavily on one or a set of reliable draft models
- ✗ Open-domain draft models has poor speculation accuracy (due to size)
- ✗ It's hard to predict query distributions and prepare specialized draft models offline to ensure speculation accuracy

Online speculative decoding (OSD): online distillation+SD



If the student proposes incorrect tokens, both the draft and target distributions are stored in a **buffer**. Once the buffer exceeds a specified threshold, the draft model is **updated** by calculating the loss between the draft and target distributions using various distance metrics.

Distillation loss

$$\ell(\boldsymbol{\theta}) = \frac{1}{n_B} \sum_{\mathbf{x}^{(i)} \in \mathcal{B}} \ell(\mathbf{x}^{(i)}, \boldsymbol{\theta}), \quad \ell(\mathbf{x}, \boldsymbol{\theta}) = D(p(\cdot|\mathbf{x}) \| q_{\boldsymbol{\theta}}(\cdot|\mathbf{x}))$$

$$\ell_{KL}(\mathbf{x}, \boldsymbol{\theta}) = D_{\text{KL}}(p(\cdot|\mathbf{x}) \| q_{\boldsymbol{\theta}}(\cdot|\mathbf{x})),$$

$$\ell_{RKL}(\mathbf{x}, \boldsymbol{\theta}) = D_{\text{KL}}(q_{\boldsymbol{\theta}}(\cdot|\mathbf{x}) \| p(\cdot|\mathbf{x})),$$

$$\ell_{JSD[\beta]}(\mathbf{x}, \boldsymbol{\theta}) = \beta D_{\text{KL}} \left(p(\cdot|\mathbf{x}) \| p_{\boldsymbol{\theta}}^{\beta}(\cdot|\mathbf{x}) \right) + (1 - \beta) D_{\text{KL}} \left(q_{\boldsymbol{\theta}}(\cdot|\mathbf{x}) \| p_{\boldsymbol{\theta}}^{\beta}(\cdot|\mathbf{x}) \right)$$

* Estimating the above objectives involves the expectation over $q_{\boldsymbol{\theta}}(\cdot|\mathbf{x})$ or $p(\cdot|\mathbf{x})$, which should be expanded **recursively**

* When sampling from $q_{\boldsymbol{\theta}}(\cdot|\mathbf{x})$, we should **differentiate through the sampling process** for unbiased gradient estimation

The algorithm

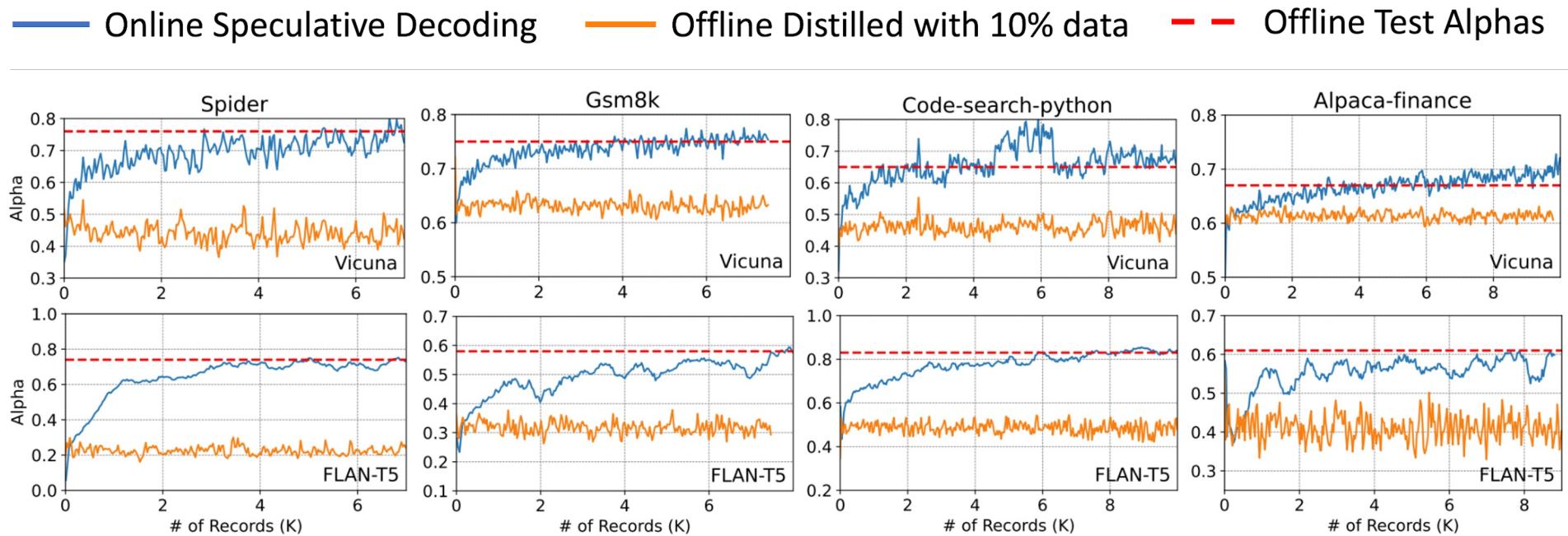
Algorithm 1 Online Speculative Decoding.

```
1: Input: Target LLM  $p(\cdot|\mathbf{x})$ , draft LLM  $q_{\theta}(\cdot|\mathbf{x})$ , warmup dataset  $\mathcal{D}$ , online data stream  $\mathcal{S}$ , guess number  $k$ ,  
temporary buffer  $\mathcal{R}$ , replay buffer  $\mathcal{Q}$ , update interval for the draft model  $I$ .  
2: Pre-train  $q_{\theta}$  to approximate  $p$  with data from  $\mathcal{D}$  by minimizing  $\ell(\mathbf{x}, \theta)$  using Equation (5);  
3:  $i \leftarrow 0$ ;  
4:  $\mathcal{Q} \leftarrow []$ ;  
5:  $cur\_len = |\mathbf{x}|$  // Total sequence length, including prompt length and tokens generated so far.  
6: while True do  
7:    $\mathcal{R} \leftarrow []$  // List of (error_index, target logits at error_index) pairs for a single request.  
8:    $\mathbf{x} \sim \mathcal{S}$ ,  $i \leftarrow i + 1$ ;  
9:   while  $\langle \text{EOS} \rangle$  not in  $\mathbf{x}$  do  
10:     $\mathbf{y} = \{y_1, \dots, y_k\} \sim q_{\theta}(\cdot|\mathbf{x})$ ;  
11:    Estimate  $\{p(y|\mathbf{x}, \mathbf{y}_{<i})\}_{i=1}^{k+1}$  in parallel;  
12:    Determine number of accepted tokens  $a$  and sample one more token, yielding  $\mathbf{y} = \{y_1, \dots, y_{a+1}\}$ ;  
13:     $cur\_len \leftarrow cur\_len + a + 1$ ;  
14:     $error\_index \leftarrow cur\_len$ ;  
15:    Append (error_index,  $p(y|\mathbf{x}, \mathbf{y}_{<a+1})$ ) to  $\mathcal{R}$ ;  
16:     $\mathbf{x} \leftarrow [\mathbf{x}, \mathbf{y}_{<a+2}]$ ;  
17:   end while  
18:   Append ( $\mathbf{x}, \mathcal{R}$ ) to  $\mathcal{Q}$ ;  
19:   if  $i \bmod I = 0$  then  
20:     Update  $q_{\theta}$  on  $\mathcal{Q}$  to minimize  $\ell(\mathbf{x}, \theta)$  analytically;  
21:      $\mathcal{Q} \leftarrow []$ ;  
22:   end if  
23: end while
```

Experimental setup

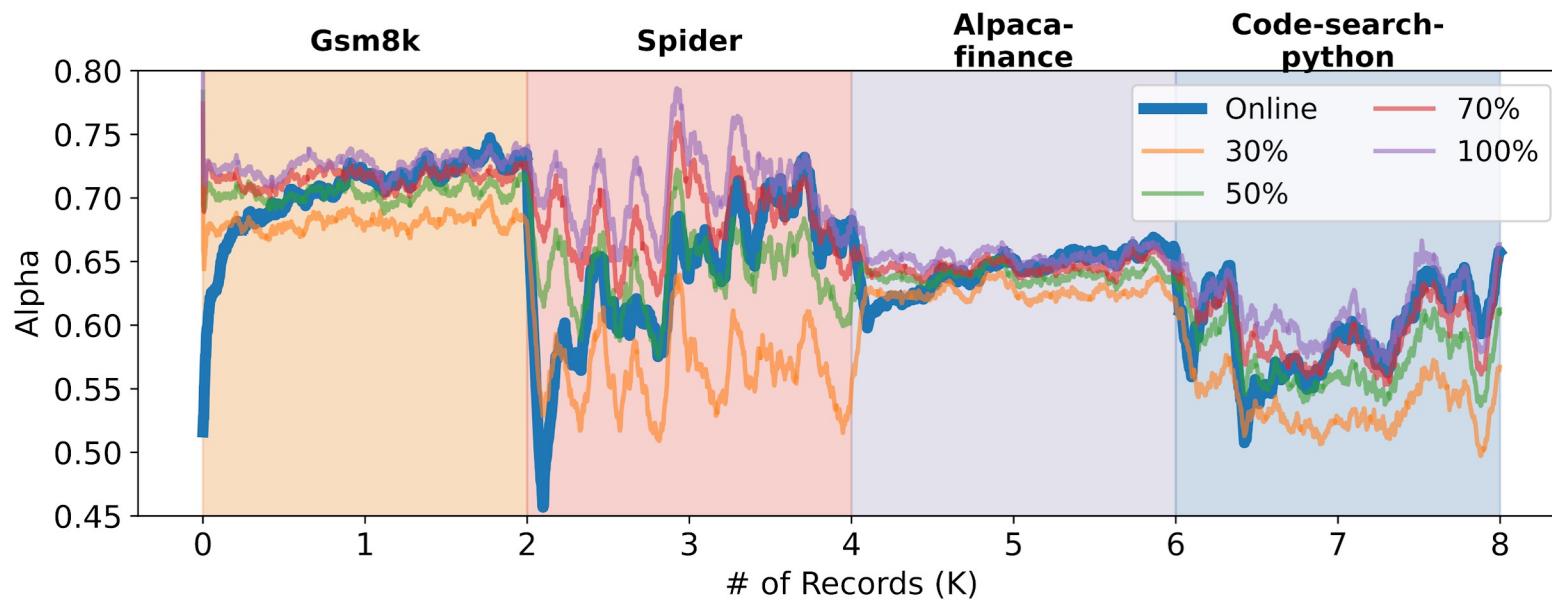
- Metric: token acceptance rate and wall-clock time
- Target model: Vicuna7B and FLAN-T5-XL (3B)
- Draft model: LLaMA-160m and T5-Small
- GPU: A100-80GB
- Datasets: Text-to-SQL (Spider), graduate school math (Gsm8k), Python code generation (Code-search-Python), and financial question answering (Alpaca-finance)
- Number of proposed tokens: 5

Does the online algorithm increase the token acceptance rate?



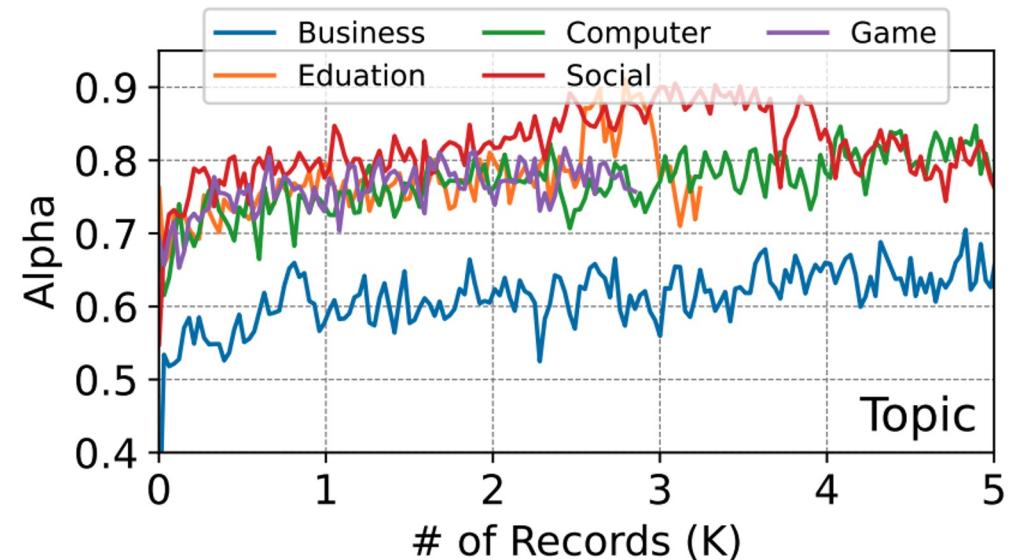
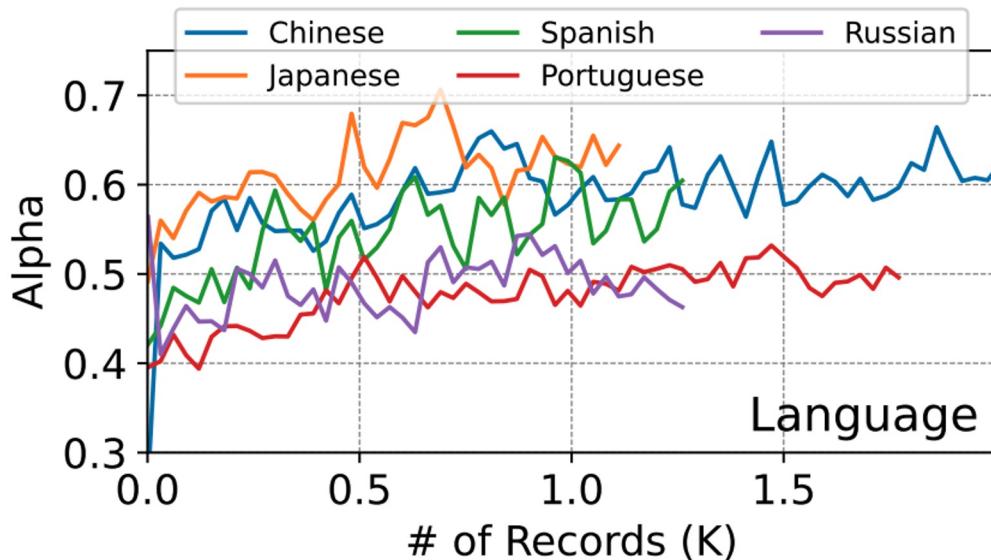
The acceptance rate rises swiftly as the draft model is exposed to more data

How quickly can OSD adapt to distribution shift



OSD's alpha value dips notably at distribution boundaries but rebounds quickly as OSD processes more data, quickly matching or even surpassing performances seen with 70% to 100% data access

OSD+routing on Arena (real LMSYS-chat conversations that span 4 months)



- One individual draft model for each language/topic
- When routing by language, OSD enhances rates by 0.1 to **0.2**
- When routing by topic, acceptance rates are above **0.6** across topics, with Social and Computer discussions peaking near **0.8**

Measured execution time/speedup and theoretical execution time/speedup

	Original	OSD, $\alpha = 0.5$	OSD, $\alpha = 0.6$	OSD, $\alpha = 0.7$	OSD, $\alpha = 0.8$	OSD, $\alpha = 0.9$
Measured time in ms/token (speedup)	51.09	39.90 (1.28 ×)	35.48 (1.44 ×)	30.96 (1.65 ×)	25.42 (2.01 ×)	19.43 (2.63 ×)
Theoretical time in ms/token (speedup)	51.09	39.00 (1.31 ×)	32.12 (1.59 ×)	26.07 (1.96 ×)	20.77 (2.46 ×)	16.38 (3.12 ×)

- Up to **2.63x speedup** over inference without speculative decoding

Measured speedup on four evaluated datasets on a single A100-80G

Dataset	Spider	Gsm8k	Alpaca-Finance	Code-Python
Measured time in ms/token (Speedup)	23.53 (2.17 ×)	27.40 (1.89 ×)	26.53 (1.92 ×)	30.12 (1.69 ×)

- TinyLLaMA-1.1B, Vicuna-33B
- Inference without speculative decoding has a token latency of 51.09 ms/token
- Up to **2.17x speedup**

Thanks